

Cultural Differences in Bias? Origin and Gender Bias in Pre-Trained German and French Word Embeddings

Mascha Kurpicz-Briki
Bern University of Applied Sciences
Biel/Bienne, Switzerland
mascha.kurpicz@bfh.ch

Abstract

Smart applications often rely on training data in form of text. If there is a bias in that training data, the decision of the applications might not be fair. Common training data has been shown to be biased towards different groups of minorities. However, there is no generic algorithm to determine the fairness of training data. One existing approach is to measure gender bias using word embeddings. Most research in this field has been dedicated to the English language. In this work, we identified that there is a bias towards gender and origin in both German and French word embeddings. In particular, we found that real-world bias and stereotypes from the 18th century are still included in today's word embeddings. Furthermore, we show that the gender bias in German has a different form from English and there is indication that bias has cultural differences that need to be considered when analyzing texts and word embeddings in different languages.

about the gender has to be made. It is therefore highly relevant to identify and mitigate gender bias in natural language processing (Sun et al., 2019).

Word embeddings are applied in several types of applications and enhance the development of machine learning and natural language processing. However, they also amplify existing social stereotypes in the human-generated training data.

Different approaches to identify and mitigate bias in word embeddings have been developed. A word embedding is a vectorial representation of a word (or phrase), trained on co-occurrences in a text corpora. Each word w is represented as a d -dimensional word vector $\vec{w} \in \mathbb{R}^d$ (Bolukbasi et al., 2016), where often $d = 300$ (Caliskan et al., 2017). In such a vector space, words with similar meaning have vectors that are close (i.e. they have a small vector distance). It has been confirmed that the vector distance can be used to represent the relationship between two words (Mikolov et al., 2013c). Using this method, problems like the following can be solved: *man is to king as woman is to x*. With simple arithmetic on vectors this problem can be solved by proposing $x = \text{queen}$ (Bolukbasi et al., 2016), because

$$\vec{man} - \vec{woman} \approx \vec{king} - \vec{queen}.$$

1 Introduction

Bias is an important topic in machine learning applications, and in particular in natural language processing. For example, it can be easily shown in automatic translation. As shown in Figure 1, when translating "She is an engineer. He is a nurse." to Turkish and then back to English, we obtain "He's an engineer. She is a nurse.". Due to the fact that in Turkish there is no difference between *he* and *she*, when translating back to English, a guess

Even if not perfectly equal to any vector in the vocabulary, the closest vector to the resultant will often be the answer to the question (Hapke et al., 2019). This is useful for different types of applications, for example word embeddings are an important source of evidence for document ranking (Nalisnick et al., 2016) (Mitra et al., 2016). However, this relationship between words can also contain problematic associations. Research demonstrated that words like *he* or *man* are associated to jobs like programmer or doctor, whereas words

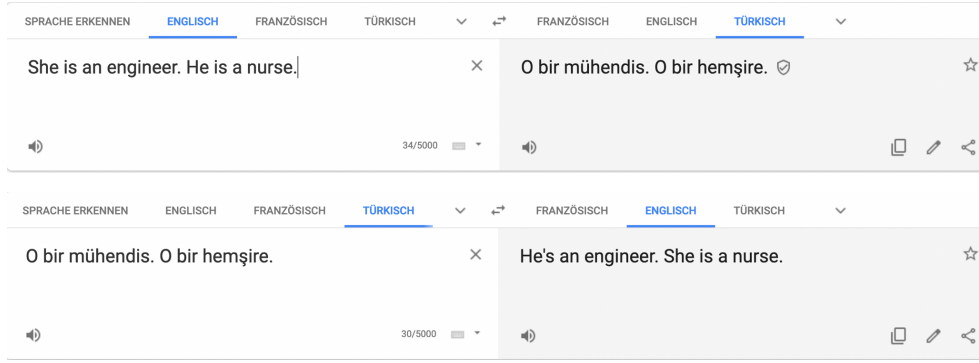


Figure 1: Example of bias in Google Translate.

like *she* or *woman* are associated to jobs like homemaker or nurse (Bolukbasi et al., 2016) (Lu et al., 2018). For example, it has been shown (Bolukbasi et al., 2016) that

$$\frac{\vec{man} - \vec{woman}}{\vec{computerprogrammer} - \vec{homemaker}} \approx$$

Human bias in psychology is often measured using Implicit Association Test (IAT) (Greenwald et al., 1998). The IAT measures differences in the response time of the human subjects, when they are asked to pair two concepts. Whenever they find these concepts similar, the response time is shorter than when they find the concepts different. Based on these results, a corresponding measure based on word embeddings instead of human subjects has been developed, called Word Embedding Association Test (WEAT) (Caliskan et al., 2017). The WEAT allows to demonstrate different types of bias in word embeddings, replacing the reaction time from IAT with word similarity (i.e. distance between word vectors). The method has been further developed and applied (e.g. (Karve et al., 2019) (May et al., 2019)), but mostly for the English language and gender bias. We apply this method to pre-trained word embeddings in German and French, and address the following research questions:

- Can known gender and origin bias found in pre-trained English word embeddings be confirmed for German and French?
- Can we identify different forms of gender bias in German word embeddings?

The paper will first discuss the related work and provide more details about the used methods. We will then describe the experimental setup. In the end, the results will be presented and discussed.

2 Related Work

2.1 Word Embeddings

Unless a domain-specific word model is required, pre-trained word vector representations are sufficient, and are easily available online as open-source (Hapke et al., 2019). In the following paragraphs we shortly describe the most common word embedding training techniques:

word2vec was first presented in 2013 (Mikolov et al., 2013b) (Mikolov et al., 2013a) (Mikolov et al., 2013c). These word embeddings provided a surprising accuracy improvement on several NLP tasks, and can be trained in two different ways (Hapke et al., 2019): with the *skip-gram* approach using a word of interest as an input, or with the *continuous bag-of-words* approach using nearby words as input.

GloVe provides another technology for generating word embeddings (Pennington et al., 2014). Whereas word2vec relies on a neural network with backpropagation, GloVe uses direct optimization.

fastText provides an improvement to word2vec (Bojanowski et al., 2017). Instead of predicting the surrounding words, it predicts the surrounding n-character grams. This results in the advantage to handle rare words much better than the original approach (Hapke et al., 2019). Pre-trained models are available in 157 languages (Grave et al., 2018).

2.2 Bias Identification in Training Data

There is a concern that artificial intelligence and smart decision making will amplify cultural stereotypes (Barocas and Selbst, 2016). Due to historical unfairness, which is represented in the training data, unfair decisions can be made in the future. Research has shown that such bias can

be identified, for example by using bayesian networks (Mancuhan and Clifton, 2014). Commonly used datasets such as Wikipedia have been proven to be biased (Wagner et al., 2015) (Wagner et al., 2016). In particular, it was also shown how dialect can lead to racial bias in common training data for hate speech detection (Sap et al., 2019).

Recent research concentrates on bias identification in word embeddings. The state-of-the-art will be presented in the next subsection.

2.3 Bias Identification in Word Embeddings

In the original WEAT paper (Caliskan et al., 2017), several different IAT results have been confirmed on pre-trained GloVe and word2vec word embeddings for the English language. Due to their experiments on off-the-shelf machine learning components, they demonstrate that cultural stereotypes have already propagated to state-of-the-art artificial intelligence applications. The WEAT has become a common method to measure bias in word embeddings, being used as a metric when developing methods to reduce bias in word embeddings (Karve et al., 2019). The authors identified different biases, in particular the following categories of gender bias: career vs. family activities, Maths vs. Arts and Science vs. Arts. Furthermore, they detected racial bias concerning African-Americans by comparing European American and African American names.

Other research proposed a framework for temporal analysis of word embeddings and observed bias changing over time and relating it to historical events (Garg et al., 2018). The approach helped to quantify stereotypes and attitudes towards women and ethnic minorities in the United States in the 20th and 21st century.

The WEAT has also been applied to word embeddings that were trained for different specific domains (Twitter, Wikipedia-based gender-balanced corpus GAP, PubMed and Google News) (Chaloner and Maldonado, 2019). The authors confirmed a statistically significant gender bias for all experiments on the Google News corpus (and for some of the experiments on the other corpora).

It has been shown that current bias mitigation methods cannot directly be applied to languages with grammatical gender such as French or Spanish (Zhou et al., 2019). However, the authors show that different types of bias can still be identified for those languages. They also present

the Modified Word Embedding Association Test (MWEAT), which is then used to evaluate the bias in the Spanish language.

The WEAT was extended to measure bias in state-of-the-art sentence encoders (May et al., 2019). The Sentence Encoder Association test (SEAT) enters the words from the WEAT experiments into sentence templates such as "This is a[n] <word>". The results suggest that recent sentence encoders exhibit less bias than previous models, but future research to further clarify this is suggested. The research focusses on English sentences only. As WEAT, SEAT can only detect presence of bias, but not its absence.

Other research (Friedman et al., 2019) identifies gender bias in word embeddings trained on Twitter data from 99 countries and 51 U.S. regions. The results are then validated against statistical gender gaps in 18 international and 5 U.S. based statistics. In this research only tweets in English were considered.

It has been explored (McCurdy and Serbetci, 2017) whether word embeddings in languages with grammatical gender show the same topical semantic bias as in English. In particular, the authors show that for German there is a positive differential association, but the WEAT shows reliable effects only for the evaluated natural gender languages English and Dutch. The training data was prepared from the OpenSubtitles corpus (Lison and Tiedemann, 2016) with translations in German, Spanish, Dutch and English.

3 Method

3.1 WEAT method

The terminology of WEAT (Caliskan et al., 2017) is borrowed from the Implicit Association Test (IAT) (Greenwald et al., 1998) from psychology. The IAT measures a person's subconscious association between concepts and therefore gives a measure for implicit bias. It is a computer-based measure, where users are asked to rapidly categorize two target concepts with an attribute. The IAT questions are based on combining possible answers to parallel non-biased questions, and therefore implicit stereotypes can be assessed. Easier pairing (i.e., shorter reaction time) is interpreted as stronger association between the concepts.

In the background, the experiment consists of two sets of *target words*, as for example (*math, algebra, ...*) and (*art, poetry, ...*).

Furthermore, two sets of *attribute words* are defined, as for example (*man, male, ...*) and (*woman, female, ...*)

In WEAT, the distance between vectors corresponds to the reaction time in IAT. As a measure of distance between the vectors, the cosine similarity between the vectors is used.

The null hypothesis is that there is no difference between the two sets of target words with regard to relative similarity to the two sets of attribute words. In other words, there is no bias between the genders regarding the target word groups.

The WEAT test can be formalized as follows (Caliskan et al., 2017): X and Y are the two sets of target words of equal size. A and B are the two sets of attribute words. $s(X, Y, A, B)$ is the test statistics.

$$s(X, Y, A, B) = \sum_{x \in X} s(x, A, B) - \sum_{y \in Y} s(y, A, B) \quad (1)$$

where

$$s(w, A, B) = \text{mean}_{a \in A} \cos(\vec{w}, \vec{a}) - \text{mean}_{b \in B} \cos(\vec{w}, \vec{b})$$

$s(w, A, B)$ measures the association of w with the attribute. $s(X, Y, A, B)$ measures the differential association of the two sets of target words with the attribute. In the equation, $\cos(\vec{a}, \vec{b})$ defines the cosine of the angle between the vectors \vec{a} and \vec{b} , which we use to measure the distance between the two vectors.

In WEAT, a permutation test is used to measure the (un)likelihood of the null hypothesis, i.e. they compute the probability that a random permutation of the attribute words would produce the observed (or greater) difference in sample means.

$\{(X_i, Y_i)\}$ denotes all the partitions of $X \cup Y$ into two sets of equal size. The one-sided p-value is then defined as (Caliskan et al., 2017):

$$Pr_i[s(X_i, Y_i, A, B) > s(X, Y, A, B)] \quad (2)$$

In our implementation, instead of the full permutation test we implemented a randomization test with 100'000 iterations, following (Chaloner and Maldonado, 2019).

The effect size is computed as Cohen's d (as for the original IAT). The effect size d is computed as (Caliskan et al., 2017)

$$\frac{\text{mean}_{x \in X} s(x, A, B) - \text{mean}_{y \in Y} s(y, A, B)}{\text{stddev}_{w \in X \cup Y} s(w, A, B)} \quad (3)$$

3.2 Experimental Setup

This section describes the different experiments we executed in our implementation of the WEAT and pre-trained word embeddings in different languages.

3.2.1 Validation: WEAT experiments

To validate our implementation, we executed selected experiments in English (WEAT 5 for origin bias and WEAT 6-8 for gender bias) from the original WEAT paper (Caliskan et al., 2017).

In a first experiment, European American and African American names are used, along with pleasant and unpleasant attributes (WEAT5-ori, detailed setup in Table 1).

We then defined the targets as male and female names and the attributes as words regarding career and family (WEAT6-ori, detailed setup in Table 2).

Another experiment considers words from maths and arts as targets, and female and male terms as attributes. Table 4 shows the exact terms of the experiment. We first executed this experiments in its original form (WEAT7-ori). We then also executed it in a reduced form (words in italic were skipped), in order to match what the German and French experiments explained in the next sections (WEAT7-mod).

We then executed an experiment that considers words from science and arts as targets, and male and female attributes. Table 4 shows the exact terms of the experiment. We first executed this experiments in its original form (WEAT8-ori). We then also executed it in a reduced form (words in italic were skipped), in order to match what the German and French experiments explained in the next sections (WEAT8-mod).

WEAT 5-7 are based on an existing Implicit Association Test (IAT) from literature (Nosek et al., 2002a), as well as WEAT 8 (Nosek et al., 2002b).

Group	WEAT5-ori	WEAT5-ger	WEAT5-fr
Group 1	Brad, Brendan, Geoffrey, Greg, Brett, Jay, Matthew, Neil, Todd, Allison, Anne, Carrie, Emily, Jill, Laurie, Kristen, Meredith, Sarah	Peter, Daniel, Hans, Thomas, Andreas, Martin, Markus, Michael, Maria, Anna, Ursula, Ruth, Monika, Elisabeth, Verena, Sandra	Jean, Daniel, Michel, Pierre, David, Philippe, Nicolas, José, Maria, Marie, Anne, Catherine, Nathalie, Ana, Isabelle, Christine
Group 2	Darnell, Hakim, Jermaine, Kareem, Jamal, Leroy, Rasheed, Tremayne, Tyrone, Aisha, Ebony, Keisha, Kenya, Latonya, Lakisha, Latoya, Tamika, Tanisha	Ladina, Fatima, Fatma, Alma, Soraya, Svetlana, Elif, Vesna, Mehmet, Mustafa, Aleksandar, Mohamed, Ibrahim, Dragan, Hasan, Mohammad	Ladina, Fatima, Fatma, Alma, Soraya, Svetlana, Elif, Vesna, Mehmet, Mustafa, Aleksandar, Mohamed, Ibrahim, Dragan, Hasan, Mohammad
Pleasant	joy, love, peace, wonderful, pleasure, friend, laughter, happy	Spass, Liebe, Frieden, wunderbar, Freude, Lachen, glücklich	joie, amour, paix, magnifique, plaisir, ami, rire, enthousiaste
Unpleasant	agony, terrible, horrible, nasty, evil, war, awful, failure	Qual, furchtbar, schrecklich, übel, böse, Krieg, scheusslich, Versagen	souffrance, terrible, horrible, désagréable, mal, guerre, abominable, défaillance

Table 1: The terms from the original WEAT 5 experiment (Caliskan et al., 2017) and our adaptations/translations to German and French.

Group	WEAT6-ori	WEAT6-ger1	WEAT6-fr1
Male names	John, Paul, Mike, Kevin, Steve, Greg, Jeff, Bill	Peter, Daniel, Hans, Thomas, Andreas, Martin, Markus, Michael	Jean, Daniel, Michel, Pierre, David, Philippe, Nicolas, José
Female names	Amy, Joan, Lisa, Sarah, Diana, Kate, Ann, Donna	Maria, Anna, Ursula, Ruth, Monika, Elisabeth, Verena, Sandra	Maria, Marie, Anne, Catherine, Nathalie, Ana, Isabelle, Christine
Career	executive, management, professional, corporation, salary, office, business, career	Führungskraft, Verwaltung, beruflich, Konzern, Gehalt, Büro, Geschäft, Werdegang	équipe, gestion, profession, société, salaire, bureau, affaires, carrière
Family	home, parents, children, family, cousins, marriage, weddings, relatives	Zuhause, Eltern, Kinder, Familie, Cousinsen, Ehe, Hochzeit, Verwandtschaft	maison, parents, enfants, famille, cousins, mariage, noces, proches

Table 2: The terms from the original WEAT 6 experiment (Caliskan et al., 2017) and our adaptations/translations to German and French for names in Switzerland.

3.2.2 Reproduction of WEAT 5-8 for German

We translated and/or adapted the experiments to execute them on German pre-trained word embeddings as described in the next paragraphs.

WEAT5-ger We reproduced the origin experiment that connected names of specific origins to pleasant or unpleasant words for German. We selected originally Swiss German names by using the 8 most common names of the German part of Switzerland for women and men respectively¹. We then selected manually a list of commonly used names in Switzerland that are of different origin from the same source. These names were chosen as representatives of names of foreign origin. A German study has shown that the origin of the name has a major impact on the success of job applications (Schneider et al., 2014). Instead of focussing on the percentage of different minorities of the population, which is complicated due to regional differences, we selected commonly

used names of different origins, based on the list of the most common names in Switzerland mentioned before. The pleasant and unpleasant terms were translated to German. Table 1 shows the exact terms of the experiment.

WEAT6-ger1 and WEAT6-ger2 We reproduced the gender experiment regarding career vs. family attributes for German. In a first experiment (WEAT6-ger1), we used the 8 most common names of the German part of Switzerland for women and men respectively². In a second experiment (WEAT6-ger2), we used the most common names of adults living in Germany³. The career and family terms were translated to German. Tables 2 and 3 show the exact terms used in the experiments.

WEAT7-ger and WEAT8-ger We reproduced the gender experiment regarding Math vs. Arts

¹Bundesamt für Statistik - Vornamen der Bevölkerung nach Jahrgang, Schweiz und Sprachgebiete, 2018

²Bundesamt für Statistik - Vornamen der Bevölkerung nach Jahrgang, Schweiz und Sprachgebiete, 2018

³<https://www.beliebte-vornamen.de/49519-erwachsene.htm>

Group	WEAT6-ori	WEAT6-ger2	WEAT6-fr2
Male names	John, Paul, Mike, Kevin, Steve, Greg, Jeff, Bill	Michael, Thomas, Andreas, Peter, Stefan, Christian, Hans, Klaus	Jean, Pierre, Michel, André, Philippe, René, Louis, Alain
Female names	Amy, Joan, Lisa, Sarah, Diana, Kate, Ann, Donna	Sabine, Susanne, Petra, Monika, Claudia, Birgit, Andrea, Stefanie	Marie, Jeanne, Françoise, Monique, Catherine, Nathalie, Isabelle, Jacqueline
Career	executive, management, professional, corporation, salary, office, business, career	Führungskraft, Verwaltung, beruflich, Konzern, Gehalt, Büro, Geschäft, Werdegang	équipe, gestion, profession, société, salaire, bureau, affaires, carrière
Family	home, parents, children, family, cousins, marriage, weddings, relatives	Zuhause, Eltern, Kinder, Familie, Cousinen, Ehe, Hochzeit, Verwandtschaft	maison, parents, enfants, famille, cousins, mariage, noces, proches

Table 3: The terms from the original WEAT 6 experiment (Caliskan et al., 2017) and our adaptations/translations to German and French for names in Germany and France.

Group	WEAT7-ori/mod	WEAT7-ger	WEAT7-fr
Math	math, algebra, geometry, calculus, equations, computation, numbers, addition	Mathematik, Algebra, Geometrie, Calculus, Gleichungen, Berechnung, Zahlen, Addition	mathématiques, algèbre, géométrie, calcul, équations, calcul, nombres, addition
Arts	poetry, art, dance, literature, novel, symphony, drama, sculpture	Poesie, Kunst, Tanz, Literatur, Roman, Symphonie, Drama, Skulptur	poésie, art, danse, littérature, roman, symphonie, drame, sculpture
Male terms	male, man, boy, brother, <i>he</i> , <i>him</i> , <i>his</i> , son	männlich, Mann, Junge, Bruder, Sohn	masculin, homme, copain, frère, fils
Female terms	female, woman, girl, sister, <i>she</i> , <i>her</i> , <i>hers</i> , daughter	weiblich, Frau, Mädchen, Schwester, Tochter	féminine, femme, copine, soeur, fille

Table 4: The terms from the original WEAT 7 experiment (Caliskan et al., 2017) and our adaptations/translations to German and French.

and Science vs. Arts for German. The pronouns in the attribute terms were skipped, because of conflicts with other terms. For example, *sie* can be *she*, but also *they*; or *sein* could be *his* but also refer to the verb *to be*. We considered NASA, Einstein and Shakespeare as internationally known and kept these words for the German experiments. Tables 4 and 5 show the exact terms of the experiment.

3.2.3 Reproduction of WEAT 6-8 for French

We translated and/or adapted the experiments to execute them on French pre-trained word embeddings as described in the next paragraphs.

WEAT5-fr We reproduced the experiment that connects names of specific origins to pleasant or unpleasant words in French. We selected originally Swiss French names by using the 8 most common names of the French part of Switzerland for women and men respectively⁴. We then selected manually a list of commonly used names in Switzerland that are of different origin from the same source (as described in the experiment WEAT5-ger). The pleasant and unpleasant terms

were translated to French. In the translation, words that have the same form for male and female (e.g. *magnifique* instead of *merveilleux*) were preferred, in order to provide consistency in the number of terms used in English and German. Table 1 shows the exact terms of the experiment.

WEAT6-fr1 and WEAT6-fr2 We reproduced the gender experiment regarding career vs. family attributes for French. To translate the female and male names, in a first experiment (WEAT6-fr1), we used the 8 most common names of the French part of Switzerland for women and men⁵. In a second experiment (WEAT6-fr2) we used the most common names in metropolitan France given between 1943 and 2019⁶. The word *executive* leads to a French word with a male and a female form. It was therefore replaced by the business related word *équipe*. Tables 2 and 3 show the exact terms of the experiments.

WEAT7-fr and WEAT8-fr As in German, pronouns were skipped. Additionally, we replaced *girl/boy* with *copain/copine* (in english:

⁴Bundesamt für Statistik - Vornamen der Bevölkerung nach Jahrgang, Schweiz und Sprachgebiete, 2018

⁵Bundesamt für Statistik - Vornamen der Bevölkerung nach Jahrgang, Schweiz und Sprachgebiete, 2018

⁶<https://tinyurl.com/tkgubf5>

Group	WEAT8-ori/mod	WEAT8-ger	WEAT8-fr
Science	science, technology, physics, chemistry, Einstein, NASA, experiment, astronomy	Wissenschaft, Technologie, Physik, Chemie, Einstein, NASA, Experiment, Astronomie	science, technologie, physique, chimie, Einstein, NASA, expérience, astronomie
Arts	poetry, art, Shakespeare, dance, literature, novel, symphony, drama	Poesie, Kunst, Shakespeare, Tanz, Literatur, Roman, Symphonie, Drama	poésie, art, Shakespeare, danse, littérature, roman, symphonie, drame
Male terms	brother, father, uncle, grandfather, son, <i>he, his, him</i>	Bruder, Vater, Onkel, Grossvater, Sohn	frère, père, oncle, grand-père, fils
Female terms	sister, mother, aunt, grandmother, daughter, <i>she, hers, her</i>	Schwester, Mutter, Tante, Grossmutter, Tochter	soeur, mère, tante, grande-mère, fille

Table 5: The terms from the original WEAT 8 experiment (Caliskan et al., 2017) and our adaptations/translations to German and French.

boyfriend/girlfriend), because the French word *fille* can be both *girl* and *daughter*. For the gender-specific adjectives we picked the male version for *masculin* and the female version for *féminine*, since we expect these words to appear more frequently. We considered NASA, Einstein and Shakespeare as internationally known and kept these words for the French experiments. Tables 4 and 5 show the exact terms of the experiments.

3.2.4 Additional Gender Stereotypes in German Word Embeddings

Based on real-world bias we defined the following two additional experiments for German:

GER-1 Study choice in Switzerland is often a matter of gender. A report about equal opportunities in Switzerland (Dubach et al., 2017) indicates that at least four out of five students are female in subjects such as special pedagogy, veterinary medicine, ethnology, educational science and psychology. On the other side, in technical studies such as mechanical engineering or computer science, only around 10-20% of the students are female. In this experiment we examine whether this bias is reflected in the word embeddings. We selected the five subjects with the highest percentage of women in 2015 (Dubach et al., 2017) (special pedagogy, veterinary medicine, ethnology, educational science, psychology). We then picked the five subjects with the lowest percentage of women in 2015 (Dubach et al., 2017) (electrical engineering, mechanical engineering, computer science, microtechnology and physics). The same male and female terms as for the WEAT7 experiment which considers the different interest of men and women in arts and maths were used for this experiment. We therefore defined target and attribute word sets as shown in Table 6.

GER-2 Studies have shown the perception of the roles of men and women in the 18th century based on dictionary entries from that time (Hausen, 1981). Based on these results, a list of words describing men and women was deduced⁷. The list is separated in different categories describing the role of women and men in the society: *Bestimmung für* (engl. *intended for*), *Aktivität/Passivität* (engl. *activity/passivity*), *Tun/Sein* (engl. *doing/being*), and their characters: *Rationalität/Emotionalität* (engl. *rationality/emotionality*), *Tugenden* (engl. *virtues*). In this study we focussed on the words indicating the characters of men and women to verify whether these stereotypes are still reflected in today’s word embeddings. We therefore selected the words from the category *Rationalität/Emotionalität* for our experiment. The category *Tugenden* was skipped due to the different number of male and female words. We therefore defined the experiment as shown in Table 7.

3.3 Data Sets: Pre-trained Word Embeddings

The validation experiments in English were executed on the same pre-trained word embeddings as in the original experiments (Caliskan et al., 2017):

- GloVe pre-trained word embeddings using the "Common Crawl" corpus (300 dimensions) with 840 billion tokens⁸
- word2vec pre-trained word embeddings using Google News (300 dimensions)⁹

The German and the French experiments were

⁷https://de.wikipedia.org/wiki/Geschlechterrolle_-_Abbildung_Polarisierung_der_Geschlechterrolle_im_18._Jahrhundert

⁸<https://nlp.stanford.edu/projects/glove/>

⁹<https://code.google.com/archive/p/word2vec/>

Group	GER-1	English
Study	Elektroingenieurwesen, Maschineningenieurwesen, Informatik, Mikrotechnik, Physik	Electrical Engineering, Mechanical Engineering, Computer Science, Microtechnology, Physics
Study	Sonderpädagogik, Veterinärmedizin, Ethnologie, Erziehungswissenschaften, Psychologie	Special Pedagogy, Veterinary Medicine, Ethnology, Educational Science, Psychology
Male terms	männlich, Mann, Junge, Bruder, Sohn	male, man, boy, brother, son
Female terms	weiblich, Frau, Mädchen, Schwester, Tochter	female, woman, girl, sister, daughter

Table 6: Experiment GER-1 verifies if existing bias in study selection appears also in German word embeddings (with English translations for better readability).

Group	GER-2	English
Character	Geist, Vernunft, Verstand, Denken, Wissen, Urteilen	Mind, Rationality, Realisation, Thinking, Knowing, Judging
Character	Gefühl, Empfinden, Empfänglichkeit, Rezeptivität, Religiosität, Verstehen	Feeling, Sentiment, Receptiveness, Religiousness, Understanding
Male terms	männlich, Mann, Junge, Bruder, Sohn	male, man, boy, brother, son
Female terms	weiblich, Frau, Mädchen, Schwester, Tochter	female, woman, girl, sister, daughter

Table 7: Experiment GER-2 verifies if existing historical bias appears also in German word embeddings (with English translations for better readability).

Experiment	p-value	Effect size d	Bias detected?
GloVe			
WEAT5-ori	$< 10^{-3}$	1.36	✓
WEAT6-ori	$< 10^{-3}$	1.8	✓
WEAT7-ori	0.058	0.94	(✓)
WEAT8-ori	0.0097	1.24	✓
WEAT7-mod	0.026	1.09	✓
WEAT8-mod	0.01	1.2	✓
word2vec			
WEAT5-ori	0.02937	0.72	✓
WEAT6-ori	$< 10^{-3}$	1.88	✓
WEAT7-ori	0.039	0.99	✓
WEAT8-ori	0.008	1.24	✓
WEAT7-mod	0.04	0.99	✓
WEAT8-mod	0.008	1.24	✓

Table 8: Results of the validation: confirming the results of the original WEAT paper (Caliskan et al., 2017) for the English language on the GloVe and word2vec dataset. We report p-values (p) and absolute value of effect size (d).

Experiment	p-value	Effect size d	Bias detected?
German			
WEAT5-ger	$< 10^{-3}$	1.134	✓
WEAT6-ger1	$< 10^{-3}$	1.62	✓
WEAT6-ger2	0.003	1.44	✓
WEAT7-ger	0.65	0.23	×
WEAT8-ger	0.83	0.11	×
GER-1	$< 10^{-3}$	1.74	✓
GER-2	0.002	1.43	✓
French			
WEAT5-fr	$< 10^{-3}$	1.29	✓
WEAT6-fr1	0.14	0.75	×
WEAT6-fr2	0.03	1.03	✓
WEAT7-fr	0.2	0.62	×
WEAT8-fr	0.53	0.32	×

Table 9: Results of the German and French experiments: translated and adapted WEAT experiments and new defined experiments. We report p-values (p) and absolute value of effect size (d).

4 Results

In this work, we consider a statistically significant bias if the p-value is below 0.05, following (Chaloner and Maldonado, 2019) and (Caliskan et al., 2017).

We confirmed the bias detected by (Caliskan et al., 2017) in the WEAT 5-8 experiments for the English language (both GloVe and word2vec datasets). Table 8 lists the detailed results.

Whereas the original WEAT5 experiment con-

executed using pre-trained fastText¹⁰ word embeddings with 300 dimensions trained on CommonCrawl and Wikipedia (Grave et al., 2018). Other word embeddings were considered, but they had either less dimensions (e.g. (Kutuzov et al., 2017)) or missing words in the vocabulary which were relevant for our experiments.

¹⁰<https://fasttext.cc/docs/en/crawl-vectors.html>

sidered European American and African American names, in our experiment common Swiss names (German and French speaking area respectively) and common names in Switzerland of different origin were considered. We were able to measure statistically significant bias based on the origin of the name, in relation to pleasant and unpleasant words, for both German and French.

In the WEAT6 experiments for German, we were able to demonstrate that there is a statistically significant gender bias for the categories family and career, for the most common names from Germany and also Switzerland. For WEAT6 in French, we could not obtain statistically significant results for Switzerland. However, the WEAT method can only detect presence of bias, but not its absence. Therefore, future research is necessary to further investigate this topic. For the most common names in France, a significant bias for the WEAT6 experiment was shown.

We could not obtain statistically significant results for the word categories math vs. arts (WEAT7) and science vs. arts (WEAT8) for German and French.

However, we identified two new sets of words in German for which we could identify a statistically significant bias. On one side, we confirmed that there is a gender bias in the word categories for different subjects of study (GER-1). On the other side, historical gender bias from the 18th century was found to be still present in today's word embeddings (GER-2).

The detailed results for the German and French experiments are listed in Table 9.

5 Discussion

We confirmed existing results for gender and origin bias in English word embeddings, and examined selected word sets for German and French word embeddings. Whereas we could partially confirm the translated (and where necessary adapted) results of the English experiments for German and French, we identified new word sets for bias in German word embeddings. The identified word sets indicate that specific regional or cultural stereotypes are included in word embeddings and therefore the bias detection may vary among different languages. Future work needs to further investigate the directions proposed in this paper and extend the word sets our work has identified.

We identified a bias towards names from different origin. We can therefore confirm that stereotypes based on names present in our society, e.g. on the labour market (Schneider et al., 2014), are also existing in word embeddings. We worked with a selection of names to get a first indication, future work must further study the differences between names encoded in word embeddings. Next to the origin, it has been shown that different prejudices such as age, the attractiveness and the intelligence of the person with the corresponding name exist (Rudolph et al., 2007), or that teachers perceive students differently, based on their names (Kube, 2009). Our results indicate that there is potential to further explore existing stereotypes and prejudices in names also in word embeddings and their implication in smart decision making.

Our results on word embeddings suggest an impact on applications using machine learning or AI. Previous studies have raised the concern that such technologies may perpetuate cultural stereotypes (Barocas and Selbst, 2016) and it has been discussed whether all implicit human biases are reflected in the statistical properties of languages (Caliskan et al., 2017). Therefore, whenever we build a system that is capable of understanding or producing natural languages (e.g. text generation, machine translation), it risks to learn the stereotypes and prejudices included in the language as well. Further research to precisely measure the different types of bias in such language models and mitigate the bias is therefore required. Future work should also identify how the observed bias in word embeddings can be related to the exact text from which they originate.

6 Conclusion

Although we partially confirmed the existing gender and origin bias also in German and French word embeddings, we showed in this research that known bias in pre-trained English word embeddings comes in a different form in German. We demonstrated that real-world bias and stereotypes from the 18th century are still included in today's word embeddings in German. Our results indicate that there are cultural differences that need to be considered in future work.

The results were obtained from publicly available pre-trained embeddings. Future work to identify and mitigate bias in word embeddings in different languages is therefore highly relevant.

References

- Solon Barocas and Andrew D Selbst. 2016. Big data’s disparate impact. *Calif. L. Rev.*, 104:671.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Advances in neural information processing systems*, pages 4349–4357.
- Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.
- Kaytlin Chaloner and Alfredo Maldonado. 2019. Measuring gender bias in word embeddings across domains and discovering new gender bias word categories. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, pages 25–32.
- Philipp Dubach, Victor Legler, Mario Morger, and Heidi Stutz. 2017. *Frauen und Männer an Schweizer Hochschulen: Indikatoren zur Chancengleichheit in Studium und wissenschaftlicher Laufbahn*. SBFI.
- Scott Friedman, Sonja Schmer-Galunder, Anthony Chen, and Jeffrey Rye. 2019. Relating word embedding gender biases to gender gaps: A cross-cultural analysis. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, pages 18–24.
- Nikhil Garg, Londa Schiebinger, Dan Jurafsky, and James Zou. 2018. Word embeddings quantify 100 years of gender and ethnic stereotypes. *Proceedings of the National Academy of Sciences*, 115(16):E3635–E3644.
- Edouard Grave, Piotr Bojanowski, Prakhar Gupta, Armand Joulin, and Tomas Mikolov. 2018. Learning word vectors for 157 languages. *arXiv preprint arXiv:1802.06893*.
- Anthony G Greenwald, Debbie E McGhee, and Jordan LK Schwartz. 1998. Measuring individual differences in implicit cognition: the implicit association test. *Journal of personality and social psychology*, 74(6):1464.
- Hannes Max Hapke, Hobson Lane, and Cole Howard. 2019. Natural language processing in action. *History of the Family in Nineteenth-and Twentieth-Century Germany*, pages 51–83.
- Saket Karve, Lyle Ungar, and João Sedoc. 2019. Conceptor debiasing of word representations evaluated on weat. *arXiv preprint arXiv:1906.05993*.
- Julia Isabell Kube. 2009. *Vornamensforschung: Fragebogenuntersuchung bei Lehrerinnen und Lehrern, ob Vorurteile bezüglich spezifischer Vornamen von Grundschulern und davon abgeleitete erwartete spezifische Persönlichkeitsmerkmale vorliegen*. Ph.D. thesis.
- Andrei Kutuzov, Murhaf Fares, Stephan Oepen, and Erik Velldal. 2017. Word vectors, reuse, and replicability: Towards a community repository of large-text resources. In *Proceedings of the 58th Conference on Simulation and Modelling*, pages 271–276. Linköping University Electronic Press.
- Pierre Lison and Jörg Tiedemann. 2016. Opensubtitles2016: Extracting large parallel corpora from movie and tv subtitles.
- Kaiji Lu, Piotr Mardziel, Fangjing Wu, Preetam Amancharla, and Anupam Datta. 2018. Gender bias in neural natural language processing. *arXiv preprint arXiv:1807.11714*.
- Koray Mancuhan and Chris Clifton. 2014. Combating discrimination using bayesian networks. *Artificial intelligence and law*, 22(2):211–238.
- Chandler May, Alex Wang, Shikha Bordia, Samuel R Bowman, and Rachel Rudinger. 2019. On measuring social biases in sentence encoders. *arXiv preprint arXiv:1903.10561*.
- Katherine McCurdy and Oguz Serbetci. 2017. Grammatical gender associations outweigh topical gender bias in crosslinguistic word embeddings. *Proceedings of WiNLP*.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013a. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013b. Distributed representations of words and phrases and their compositional-ity. In *Advances in neural information processing systems*, pages 3111–3119.
- Tomáš Mikolov, Wen-tau Yih, and Geoffrey Zweig. 2013c. Linguistic regularities in continuous space word representations. In *Proceedings of the 2013 conference of the north american chapter of the association for computational linguistics: Human language technologies*, pages 746–751.
- Bhaskar Mitra, Eric Nalisnick, Nick Craswell, and Rich Caruana. 2016. A dual embedding space model for document ranking. *arXiv preprint arXiv:1602.01137*.

- Eric Nalisnick, Bhaskar Mitra, Nick Craswell, and Rich Caruana. 2016. Improving document ranking with dual word embeddings. In *Proceedings of the 25th International Conference Companion on World Wide Web*, pages 83–84.
- Brian A Nosek, Mahzarin R Banaji, and Anthony G Greenwald. 2002a. Harvesting implicit group attitudes and beliefs from a demonstration web site. *Group Dynamics: Theory, Research, and Practice*, 6(1):101.
- Brian A Nosek, Mahzarin R Banaji, and Anthony G Greenwald. 2002b. Math= male, me= female, therefore math \neq me. *Journal of personality and social psychology*, 83(1):44.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- Udo Rudolph, Robert Böhm, and Michaela Lummer. 2007. Ein vorname sagt mehr als 1000 worte. *Zeitschrift für Sozialpsychologie*, 38(1):17–31.
- Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A Smith. 2019. The risk of racial bias in hate speech detection. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1668–1678.
- Jan Schneider, Ruta Yemane, and Martin Weinmann. 2014. *Diskriminierung am Ausbildungsmarkt: Ausmaß, Ursachen und Handlungsperspektiven*. Sachverständigenrat deutscher Stiftungen für Integration und Migration GmbH . . .
- Tony Sun, Andrew Gaut, Shirlyn Tang, Yuxin Huang, Mai ElSherief, Jieyu Zhao, Diba Mirza, Elizabeth Belding, Kai-Wei Chang, and William Yang Wang. 2019. Mitigating gender bias in natural language processing: Literature review. *arXiv preprint arXiv:1906.08976*.
- Claudia Wagner, David Garcia, Mohsen Jadidi, and Markus Strohmaier. 2015. It’s a man’s wikipedia? assessing gender inequality in an online encyclopedia. In *Ninth international AAAI conference on web and social media*.
- Claudia Wagner, Eduardo Graells-Garrido, David Garcia, and Filippo Menczer. 2016. Women through the glass ceiling: gender asymmetries in wikipedia. *EPJ Data Science*, 5(1):5.
- Pei Zhou, Weijia Shi, Jieyu Zhao, Kuan-Hao Huang, Muhao Chen, Ryan Cotterell, and Kai-Wei Chang. 2019. Examining gender bias in languages with grammatical gender. *arXiv preprint arXiv:1909.02224*.